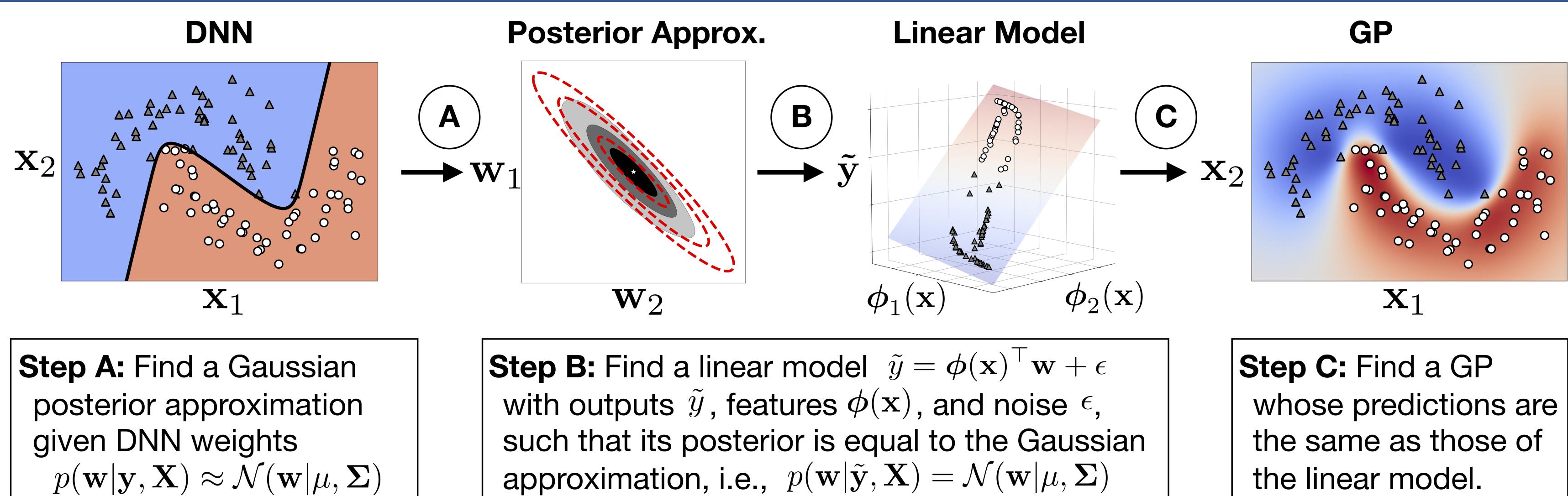


Overview

Summary:

- There is limited work and understanding on the connection between training algorithms for deep learning and Gaussian processes
- We want to relate solutions and iterations of a deep-learning algorithm to GP inference
- We show kernels obtained on real datasets and demonstrate the use of the GP marginal likelihood to tune hyperparameters of DNNs



Step A: Posterior approximation for DNNs

Deep Learning

- Training set:** $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, input $\mathbf{x}_i \in \mathbb{R}^D$, output $\mathbf{y}_i \in \mathbb{R}^K$
- Neural network:** $\mathbf{f}_\mathbf{w}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^K$, parameter vector $\mathbf{w} \in \mathbb{R}^P$
- Loss function:** $\ell(\mathbf{y}, \mathbf{f}_\mathbf{w}(\mathbf{x}))$ twice diff. and strictly convex in $\mathbf{f}_\mathbf{w}(\mathbf{x})$

$$\mathbf{w}_* = \arg \min_{\mathbf{w}} \bar{\ell}(\mathcal{D}, \mathbf{w}) := \sum_{i=1}^N \ell(\mathbf{y}_i, \mathbf{f}_\mathbf{w}(\mathbf{x}_i)) + \frac{1}{2} \delta \mathbf{w}^\top \mathbf{w},$$

We define the Jacobian $\mathbf{J}_\mathbf{w}(\mathbf{x}) := \nabla_\mathbf{w} \mathbf{f}_\mathbf{w}(\mathbf{x})^\top$, the loss residual $\mathbf{r}_\mathbf{w}(\mathbf{x}, \mathbf{y}) := \nabla_\mathbf{w} \ell(\mathbf{y}, \mathbf{f}_\mathbf{w}(\mathbf{x}))$ and the loss Hessian $\Lambda_\mathbf{w}(\mathbf{x}, \mathbf{y}) := \nabla_\mathbf{w}^2 \ell(\mathbf{y}, \mathbf{f}_\mathbf{w}(\mathbf{x}))$.

Approximations to the Posterior

Laplace GGN approximation $q_{\text{Lap}}(\mathbf{w}) := \mathcal{N}(\mathbf{w}|\mathbf{w}_*, \tilde{\Sigma})$ where

$$\tilde{\Sigma}^{-1} = \sum_{i=1}^N \mathbf{J}_{\mathbf{w}_*}(\mathbf{x}_i)^\top \Lambda_{\mathbf{w}_*}(\mathbf{x}_i, \mathbf{y}_i) \mathbf{J}_{\mathbf{w}_*}(\mathbf{x}_i) + \delta \mathbf{I}_P$$

Online GGN iterates $q_t(\mathbf{w}) := \mathcal{N}(\mathbf{w}|\mu_t, \Sigma_t)$ with $\mu := \mu_t$ below

$$\begin{aligned} \mu_{t+1} &= \mu - \beta_t \Sigma_{t+1} \left[\sum_{i=1}^N \mathbf{J}_\mu(\mathbf{x}_i)^\top \mathbf{r}_\mu(\mathbf{x}_i, \mathbf{y}_i) + \delta \mu \right], \\ \Sigma_{t+1}^{-1} &= (1 - \beta_t) \Sigma_t^{-1} + \beta_t \left[\sum_{i=1}^N [\mathbf{J}_\mu(\mathbf{x}_i)^\top \Lambda_\mu(\mathbf{x}_i, \mathbf{y}_i) \mathbf{J}_\mu(\mathbf{x}_i)] + \delta \mathbf{I}_P \right]. \end{aligned}$$

Step B: DNN Inference to Linear Model

Define a linear model with prior $\mathbf{w} \sim \mathcal{N}(0, \delta^{-1} \mathbf{I}_P)$, noise $\epsilon \sim \mathcal{N}(0, (\Lambda_{\mathbf{w}_*}(\mathbf{x}, \mathbf{y}))^{-1})$, and dataset as shown below:

Model: $\tilde{\mathbf{y}} = \mathbf{J}_{\mathbf{w}_*}(\mathbf{x}) \mathbf{w} + \epsilon$,

Data $\tilde{\mathcal{D}}$: $\tilde{\mathbf{y}}_i := \mathbf{J}_{\mathbf{w}_*}(\mathbf{x}_i) \mathbf{w}_* - (\Lambda_{\mathbf{w}_*}(\mathbf{x}_i, \mathbf{y}_i))^{-1} \mathbf{r}_{\mathbf{w}_*}(\mathbf{x}_i, \mathbf{y}_i)$.

The Laplace GGN posterior approximation of the deep network is equivalent to the exact posterior distribution of this linear model, i.e.

$$q_{\text{Lap}}(\mathbf{w}) = p(\mathbf{w}|\tilde{\mathcal{D}})$$

Define a linear model with $\mathbf{w} \sim \mathcal{N}(\mathbf{m}_t, \mathbf{S}_t)$, $\epsilon \sim \mathcal{N}(0, (\beta_t \Lambda_\mu(\mathbf{x}, \mathbf{y}))^{-1})$ where $\mu := \mu_t$, $\mathbf{S}_t := ((1 - \beta_t) \Sigma_t^{-1} + \beta_t \delta \mathbf{I}_P)^{-1}$, $\mathbf{m}_t := (1 - \beta_t) \mathbf{S}_t \Sigma_t^{-1} \mu_t$:

Model: $\tilde{\mathbf{y}}_t = \mathbf{J}_\mu(\mathbf{x}) \mathbf{w} + \epsilon$

Data $\tilde{\mathcal{D}}_t$: $\tilde{\mathbf{y}}_{t,i} := \mathbf{J}_\mu(\mathbf{x}_i) \mu - (\Lambda_\mu(\mathbf{x}_i, \mathbf{y}_i))^{-1} \mathbf{r}_\mu(\mathbf{x}_i, \mathbf{y}_i)$.

The Online GGN updates to the posterior approximation of the deep network are equivalent to exact inference in this linear model, i.e.

$$q_{t+1}(\mathbf{w}) = p(\mathbf{w}|\tilde{\mathcal{D}}_t)$$

Step C: From Linear Model to GP

Using the equivalence between the weight-space view and the function-space view, we get a GP regression model whose posterior predictive distribution $p(f(\mathbf{x}_*)|\mathbf{x}_*, \tilde{\mathcal{D}})$ is equal to that of the linear model.

For the Laplace GGN posterior approximation:

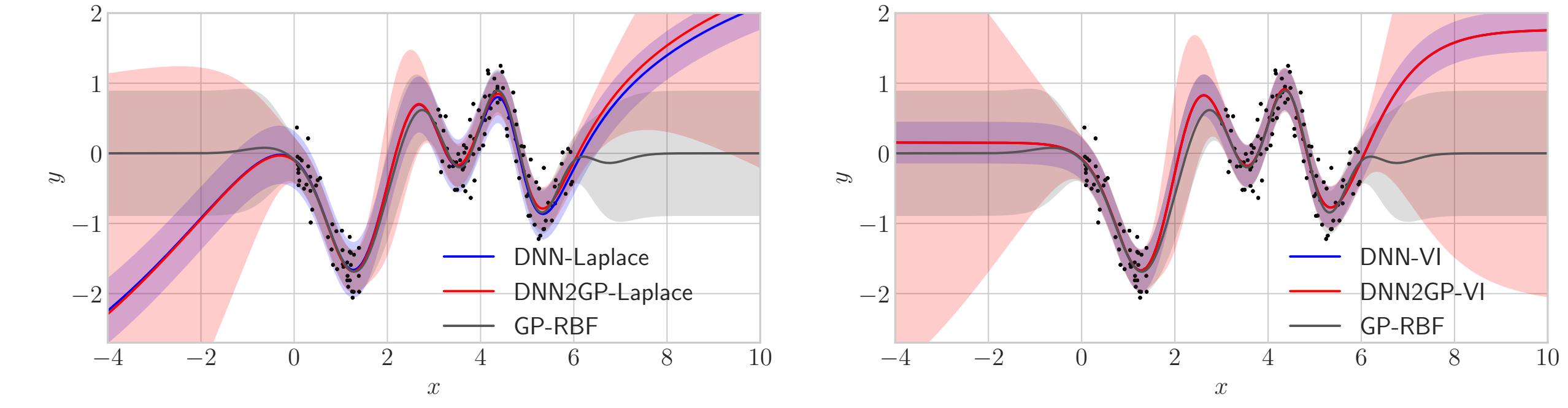
$$\tilde{\mathbf{y}} = \mathbf{f}(\mathbf{x}) + \epsilon, \text{ with } \mathbf{f}(\mathbf{x}) \sim \mathcal{GP}(0, \delta^{-1} \mathbf{J}_*(\mathbf{x}) \mathbf{J}_*(\mathbf{x}')^\top).$$

For the Online GGN updates:

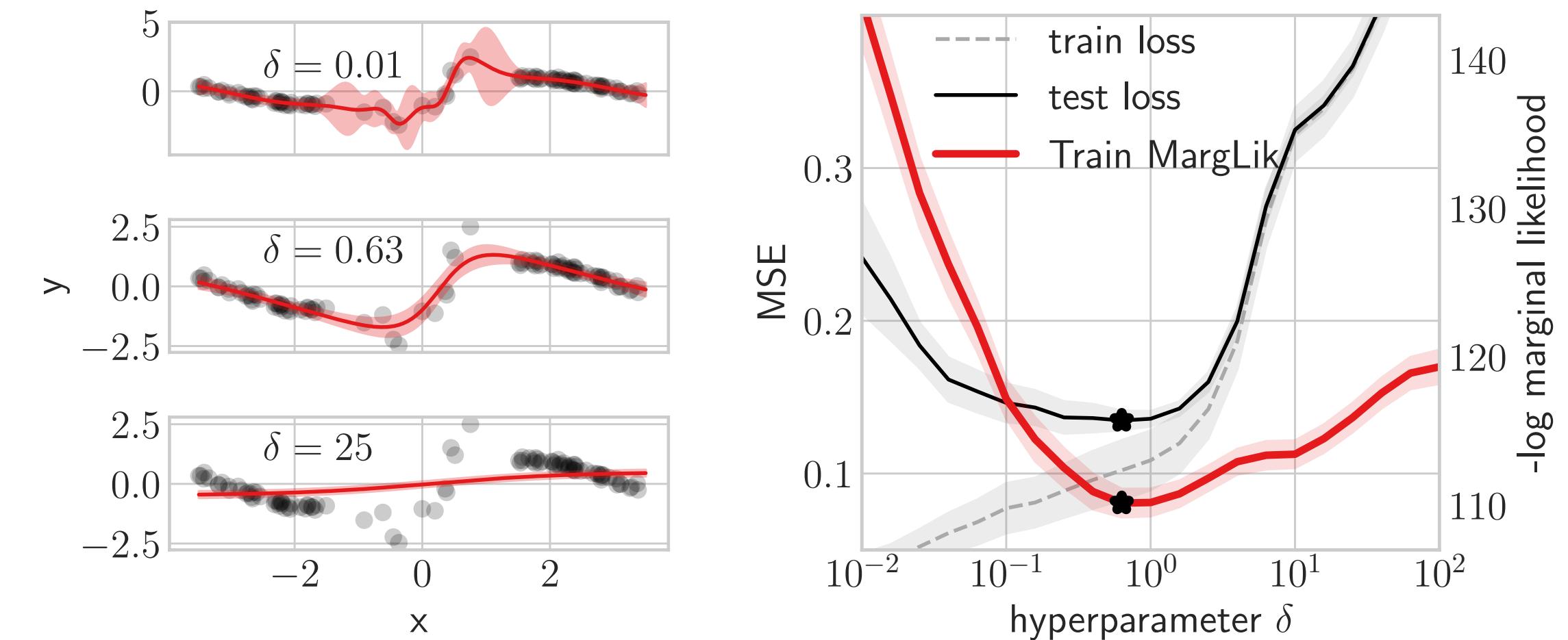
$$\tilde{\mathbf{y}}_t = \mathbf{f}_t(\mathbf{x}) + \epsilon, \text{ with } \mathbf{f}_t(\mathbf{x}) \sim \mathcal{GP}(\mathbf{J}_\mu(\mathbf{x}) \mathbf{m}_t, \mathbf{J}_\mu(\mathbf{x}) \mathbf{S}_t \mathbf{J}_\mu(\mathbf{x}')^\top).$$

Applications

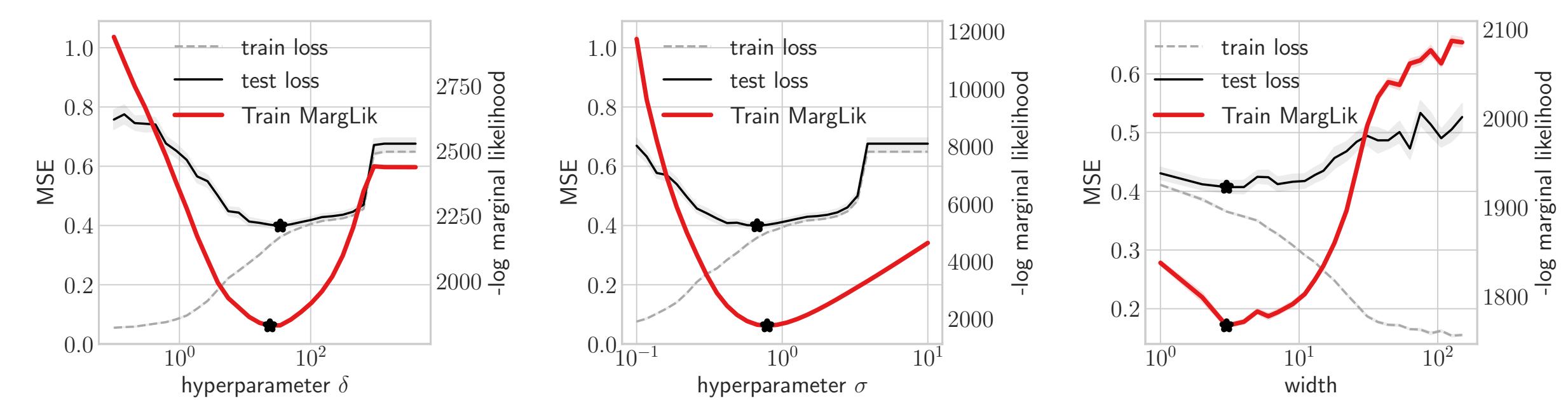
Uncertainty Estimation: by converting a neural network to a GP, we can analytically obtain high quality predictive uncertainty estimates. Below, predictive uncertainty of DNN2GP for Laplace (left) and VI (right) diagonal Gaussian approximation on a toy regression dataset:



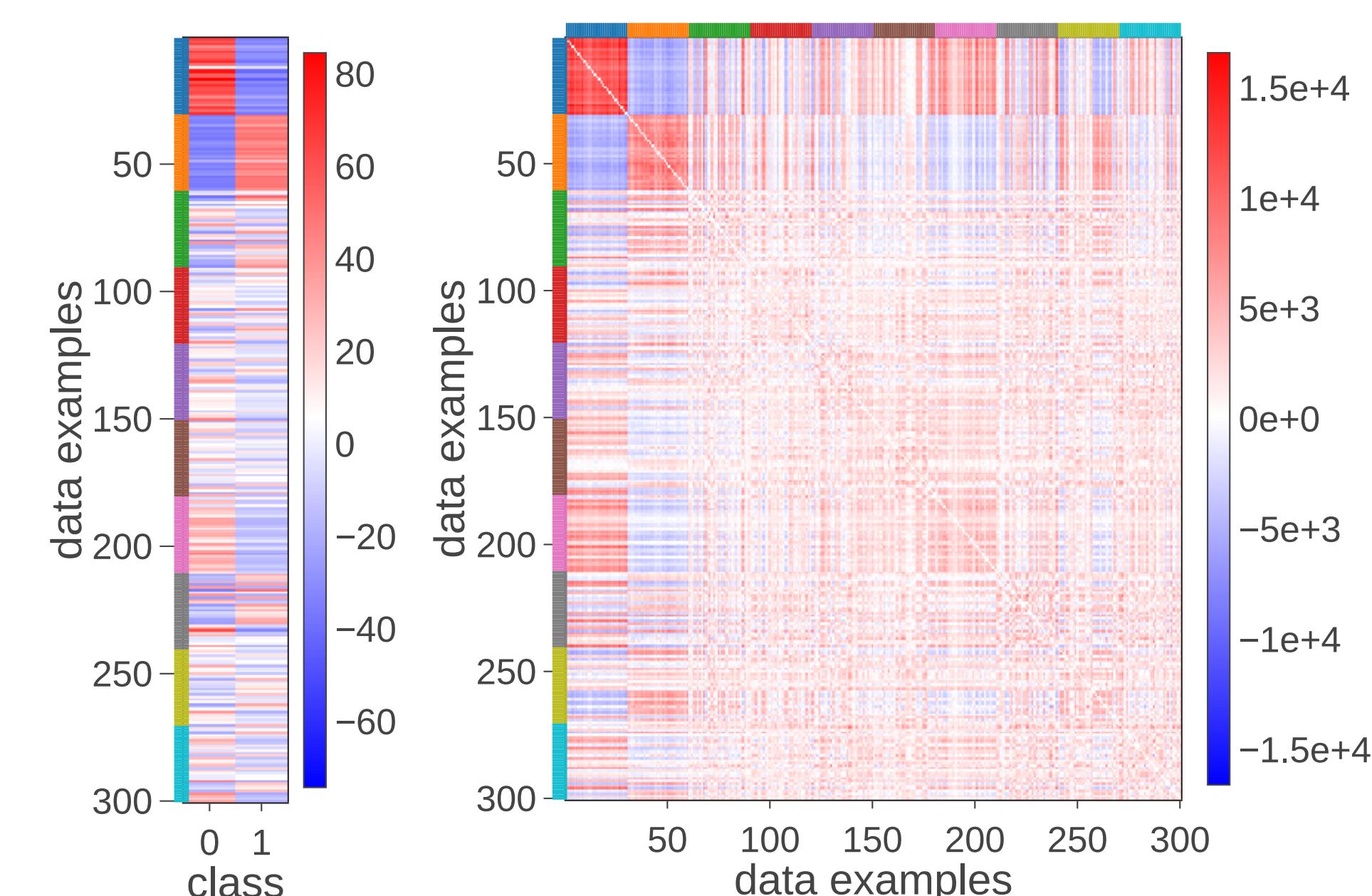
Model Selection via GP Marginal Likelihood: marginal likelihood of the GP allows Bayesian model selection for deep learning using only training data. For example, selection of δ on sinusoidal toy data:



On UCI Wine data for δ , observation noise variance σ^2 , network width:



Generalization and Interpretability: posterior mean and GP kernel on 2 out of 10 MNIST classes. We clearly see that the *in-class* digits are assigned higher posterior mean and the correlation learned is realistic.



Possible future applications: online hyperparameter tuning, coresets selection, continual learning, active learning, etc.

References

- [1] Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. SIAM Review, 2018.
- [2] Arthur Jacot, Gabriel Franck, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. NIPS, 2018.
- [3] Mohammad Emtiyaz Khan, Didrik Nielsen, Voot Tangkaratt, Wu Lin, Yarin Gal, and Akash Srivastava. Fast and Scalable Bayesian Deep Learning by Weight-Perturbation in Adam. ICML, 2018.